WHAT IS CLAIMED IS:

5          1. A document classification system for
classifying a document based on contents of the document
of which contents contains a plurality of items, said
document classification system comprising:

10          inputting means for inputting document data
corresponding to the document data;

            designating means for designating at least one
of the items contained in the document input by said
inputting means;

15          converting means for converting the document
data into converted data so that the converted data
contains only data corresponding to the item designated
by said designating means; and

            classifying means for classifying the document
20  by using the converted data produced by said converting
means.

25

2.   The document classification system as claimed in claim 1, wherein said classifying means includes document vector producing means for producing a feature vector representing a feature of the converted

5   data so as to classify the document in accordance with the feature vector produced by said document vector producing means.

10

3.   The document classification system as claimed in claim 1, wherein said converting means includes separation sign inserting means for inserting a

15   predetermined sign between sets of data corresponding to the items so as to facilitate separation of each data corresponding to each item in the converted data.

20

4.   A document classification method for classifying a document based on contents of the document of which contents contains a plurality of items, said

25   document classification method comprising the steps of:

inputting document data corresponding to the document data;

designating at least one of the items contained in the document input in the inputting step;

5    converting the document data into converted data so that the converted data contains only data corresponding to the item designated in the designating step; and

classifying the document by using the
10   converted data produced in the converting step.

15        5.   The document classification method as claimed in claim 4, wherein the classifying step includes the step of producing a feature vector representing a feature of the converted data so as to classify the document in accordance with the feature
20   vector.

25

6.  The document classification system as
claimed in claim 4, wherein the converting step includes
the step of inserting a predetermined sign between sets
of data corresponding to the items so as to facilitate

5    separation of each data corresponding to each item in
the converted data.

10

7.  A processor readable medium storing
program code causing a computer to classify a document
based on contents of the document of which contents
contains a plurality of items, comprising:

15           first program code means for inputting
document data corresponding to the document data;
             second program code means for designating at
least one of the items contained in the document;
             third program code means for converting the

20   document data into converted data so that the converted
data contains only data corresponding to the item
designated by the second program code means; and
             fourth program code means for classifying the
document by using the converted data produced by the

25   third program code means.

8. The processor readable medium as claimed in claim 7, wherein the fourth program code means includes fifth program code means for producing a feature vector representing a feature of the converted

5    data so as to classify the document in accordance with the feature vector.

10

9. The processor readable medium as claimed in claim 7, wherein the third program code means includes sixth program code means for inserting a predetermined sign between sets of data corresponding to

15    the items so as to facilitate separation of each data corresponding to each item in the converted data.

20

10. A document classification system for classifying a document according to contents of the document, said document classification system comprising:

25    input means for inputting document data of the

document;

analyzing means for analyzing the document
data so as to obtain analysis information;

vector producing means for producing a
document feature vector with respect to the document
data based on the analysis information;

transforming function calculating means for
calculating a representation transforming function used
for projecting the document feature vector onto a space
in which similarity between the document feature vectors
is reflected;

vector transforming means for transforming the
document feature vector by using the representation
transforming function;

classification means for classifying the
document based on similarity between the document
feature vectors transformed by the vector transforming
means; and

classification result storing means for
storing a result of classification performed by the
classification means.

11. The document classification system as claimed in claim 10, further comprising inner product calculating means for calculating an inner product between the document feature vectors, wherein said

5    representation transforming function calculating means calculates the representation transforming function by using the inner product.

10

12. The document classification system as claimed in claim 11, further comprising document similarity information setting means for setting

15    document similarity setting information including data representing an author of the document and a date of production of the document, wherein said representation transforming function calculating means calculates the representation transforming function by using the inner

20    product and the document similarity information.

25

13. The document classification system as claimed in claim 10, further comprising:

vector storing means for storing the document feature vector produced by said vector producing means; and

5

transforming function storing means for storing the representation transforming function calculated by said representation transforming function calculating means.

10

14. The document classification system as claimed in claim 10, further comprising vector

15 correcting means for correcting the document feature vector before the document feature vector is transformed by said vector transforming means, a correction being performed by processing one of the document feature

20 vector and a feature dimension constituting the document feature vector in accordance with a rule established by characteristics of words extracted by said analyzing means.

25

15. The document classification system as claimed in claim 14, further comprising transforming function correcting means for correcting the representation transforming function calculated by said

5    transforming function calculating means when the feature dimension is changed due to a correction of the document feature vector by said vector correcting means so that the document feature vector is transformed by said vector transforming means in accordance with the changed

10   feature dimension.



15        16. The document classification system as claimed in claim 10, further comprising:

transforming function correction instructing means for sending an instruction regarding a process to be applied on a feature dimension of the representation

20   transforming function; and

transforming function correcting means for correcting the representation transforming function based on a content of the instruction sent from said transforming function correction instructing means.

25

17. The document classification system as claimed in claim 16, wherein the process indicated in the content of the instruction is performed by using data of an arbitrary document vector.

5

18. The document classification system as claimed in claim 16, wherein the process indicated in the content of the instruction is performed by using the document feature vectors.

10

15

19. The document classification system as claimed in claim 16, wherein the process indicated in the content of the instruction is performed by using the analysis information obtained by said analyzing means.

20

25

20. The document classification system as claimed in claim 16, wherein the process indicated in the content of the instruction is performed by using the result of classification stored in said classification-
5   result storing means.

10           21. The document classification system as claimed in claim 10, further comprising:

an initial cluster centroid designating means for designating an initial cluster centroid; and

initial cluster centroid registering means for
15  registering the initial cluster centroid designated by said initial cluster centroid designating means,

wherein said classification means classifies the document in accordance with the initial cluster centroid registered by said initial cluster centroid
20  registering means.

22. The document classification system as
25  claimed in claim 21, wherein the initial cluster

centroid designated by said initial cluster centroid
designating means is arbitrary document vector data.

5

      23.   The document classification system as
claimed in claim 21, wherein the initial cluster
centroid designated by said initial cluster centroid
10   designating means is the document feature vector.

15       24.   The document classification system as
claimed in claim 21, wherein the initial cluster
centroid designated by said initial cluster centroid
designating means is the analysis information obtained
by said analyzing means.

20

      25.   The document classification system as
25   claimed in claim 21, wherein the initial cluster

centroid designated by said initial cluster centroid designating means is the result of classification stored by said classification-result storing means.

5

26.  A document classification method for classifying a document according to contents of the

10  document, said document classification method comprising the steps of:

inputting document data of the document;

analyzing the document data so as to obtain analysis information;

15  producing a document feature vector with respect to the document data based on the analysis information;

calculating a representation transforming function used for projecting the document feature vector

20  onto a space in which similarity between the document feature vectors is reflected;

transforming the document feature vector by using the representation transforming function;

classifying the document based on similarity

25  between the document feature vectors transformed in the

step of transforming; and

storing a result of classification performed in the step of classifying.

5

27. The document classification method as claimed in claim 26, further comprising the step of

10 calculating an inner product between the document feature vectors, wherein the representation transforming function is calculated by using the inner product.

15

28. The document classification method as claimed in claim 27, further comprising the step of setting document similarity setting information

20 including data representing an author of the document and a date of production of the document, wherein the representation transforming function is calculated by using the inner product and the document similarity information.

25

29. The document classification method as claimed in claim 26, further comprising the steps of:

storing the document feature vector produced in the step of producing said document feature vector;

5 and

storing the representation transforming function calculated in the step of calculating said representation transforming function.

10

30. The document classification method as claimed in claim 26, further comprising the step of

15 correcting the document feature vector before the document feature vector is transformed in the step of transforming, a correction being performed by processing one of the document feature vector and a feature dimension constituting the document feature vector in

20 accordance with a rule established by characteristics of words extracted in the step of analyzing.

25

31. The document classification method as claimed in claim 30, further comprising the step of correcting the representation transforming function calculated in the step of calculating when the feature

5　dimension is changed due to a correction of the document feature vector in the step of correcting so that the document feature vector is transformed in the step of transforming in accordance with the changed feature dimension.

10

32. The document classification method as

15　claimed in claim 26, further comprising the steps of:

sending an instruction regarding a process to be applied on a feature dimension of the representation transforming function; and

correcting the representation transforming

20　function based on a content of the instruction sent in the step of sending.

25

33. The document classification method as claimed in claim 32, wherein the process indicated in the content of the instruction is performed by using data of an arbitrary document vector.

5

34. The document classification method as

10 claimed in claim 32, wherein the process indicated in the content of the instruction is performed by using the document feature vectors.

15

35. The document classification method as claimed in claim 32, wherein the process indicated in the content of the instruction is performed by using

20 analysis information obtained by said analyzing means.

25

36. The document classification method as claimed in claim 32, wherein the process indicated in the content of the instruction is performed by using the result of classification stored in said classification-

5 result storing means.

10 37. The document classification method as claimed in claim 26, further comprising the steps of:

designating an initial cluster centroid; and

registering the initial cluster centroid designated in the step of designating,

15 wherein the document is classified in accordance with the initial cluster centroid registered in the step of registering.

20

38. The document classification method as claimed in claim 37, wherein the initial cluster centroid designated in the step of designating is

25 arbitrary document vector data.

39. The document classification method as
claimed in claim 37, wherein the initial cluster
centroid designated in the step of designating is the
document feature vector.

5

40. The document classification method as
10  claimed in claim 37, wherein the initial cluster
centroid designated in the step of designating is the
analysis information obtained in the step of analyzing.

15

41. The document classification method as
claimed in claim 37, wherein the initial cluster
centroid designated in the step of designating is the
20  result of classification stored in the step of storing.

25

42. A processor readable medium storing program code causing a computer to classify a document according to contents of the document, comprising:

first program code means for inputting

5    document data of the document;

second program code means for analyzing the document data so as to obtain analysis information;

third program code means for producing a document feature vector with respect to the document

10    data based on the analysis information;

fourth program code means for calculating a representation transforming function used for projecting the document feature vector onto a space in which similarity between the document feature vectors is

15    reflected;

fifth program code means for transforming the document feature vector by using the representation transforming function;

sixth program code means for classifying the

20    document based on similarity between the document feature vectors transformed by the fifth program code means; and

seventh program code means for storing a result of classification performed by the classification

25    means.

43.  The processor readable medium as claimed
in claim 42, further comprising eighth program code
means for calculating an inner product between the
document feature vectors, wherein the representation

5  transforming function is calculated by using the inner
product.

10

44.  The processor readable medium as claimed
in claim 43, further comprising ninth program code means
for setting document similarity setting information
including data representing an author of the document

15  and a date of production of the document, wherein the
representation transforming function is calculated by
using the inner product and the document similarity
information.

20

45.  The processor readable medium as claimed
in claim 42, further comprising:

25  tenth program code means for storing the

document feature vector produced by the third program

code means; and

eleventh program code means for storing the

representation transforming function calculated by the

5    fourth program code means.


10          46.   The processor readable medium as claimed

in claim 42, further comprising twelfth program code

means for correcting the document feature vector before

the document feature vector is transformed by the fifth

program code means, a correction being performed by

15   processing one of the document feature vector and a

feature dimension constituting the document feature

vector in accordance with a rule established by

characteristics of words extracted by the second program

code means.

20


          47.   The processor readable medium as claimed

in claim 46, further comprising thirteenth program code

25   means for correcting the representation transforming

function calculated by the fourth program code means

when the feature dimension is changed due to a

correction of the document feature vector by the twelfth

program code means so that the document feature vector

5   is transformed by the fifth program code means in

accordance with the changed feature dimension.

10

48.   The processor readable medium as claimed

in claim 42, further comprising:

fourteenth program code means for sending an

instruction regarding a process to be applied on a

15   feature dimension of the representation transforming

function; and

fifteenth program code means for correcting

the representation transforming function based on a

content of the instruction sent by the fourteenth

20   program code means.

25

49.   The processor readable medium as claimed in claim 42, further comprising:

sixteenth program code means for designating an initial cluster centroid; and

5      seventeenth program code means for registering the initial cluster centroid designated by the sixteenth program code means,

wherein the document is classified in accordance with the initial cluster centroid registered

10    by the seventeenth program code means.

15

20

25